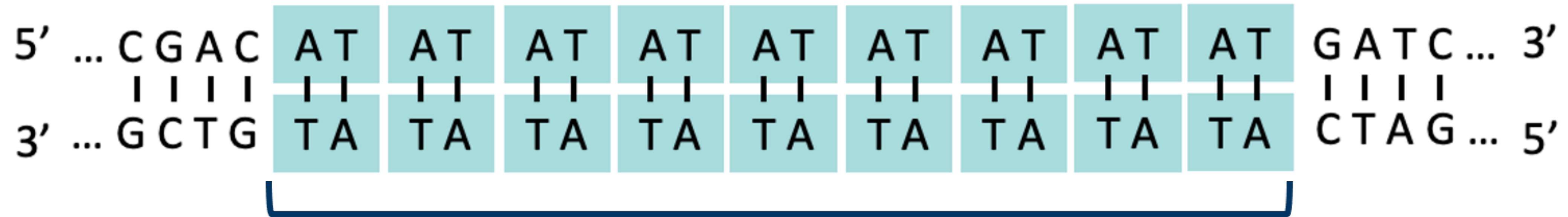


Analysing short tandem repeats in genomic sequences

What, why, how?

Short tandem repeats (a.k.a. microsatellites)

Repetitions of 1-6 nt DNA motifs



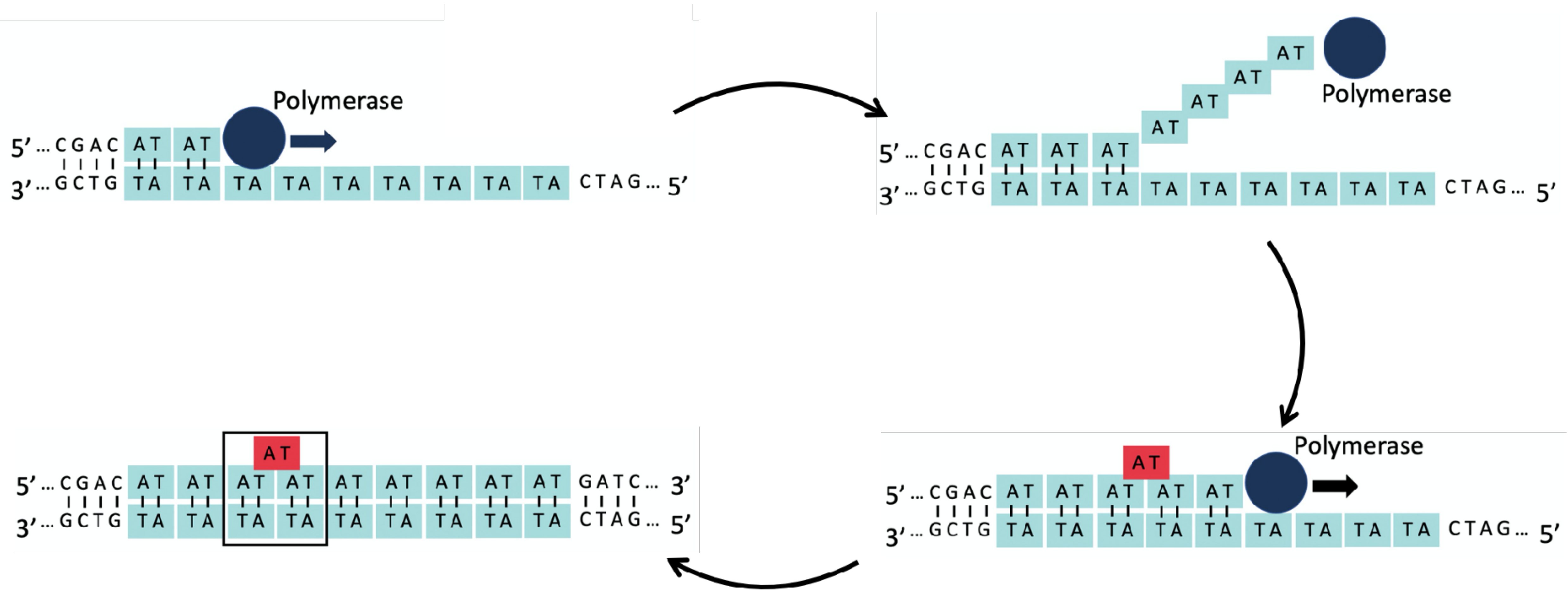
Motif: the repeating unit ('AT' in this case)

Motif length: length of the repeating unit

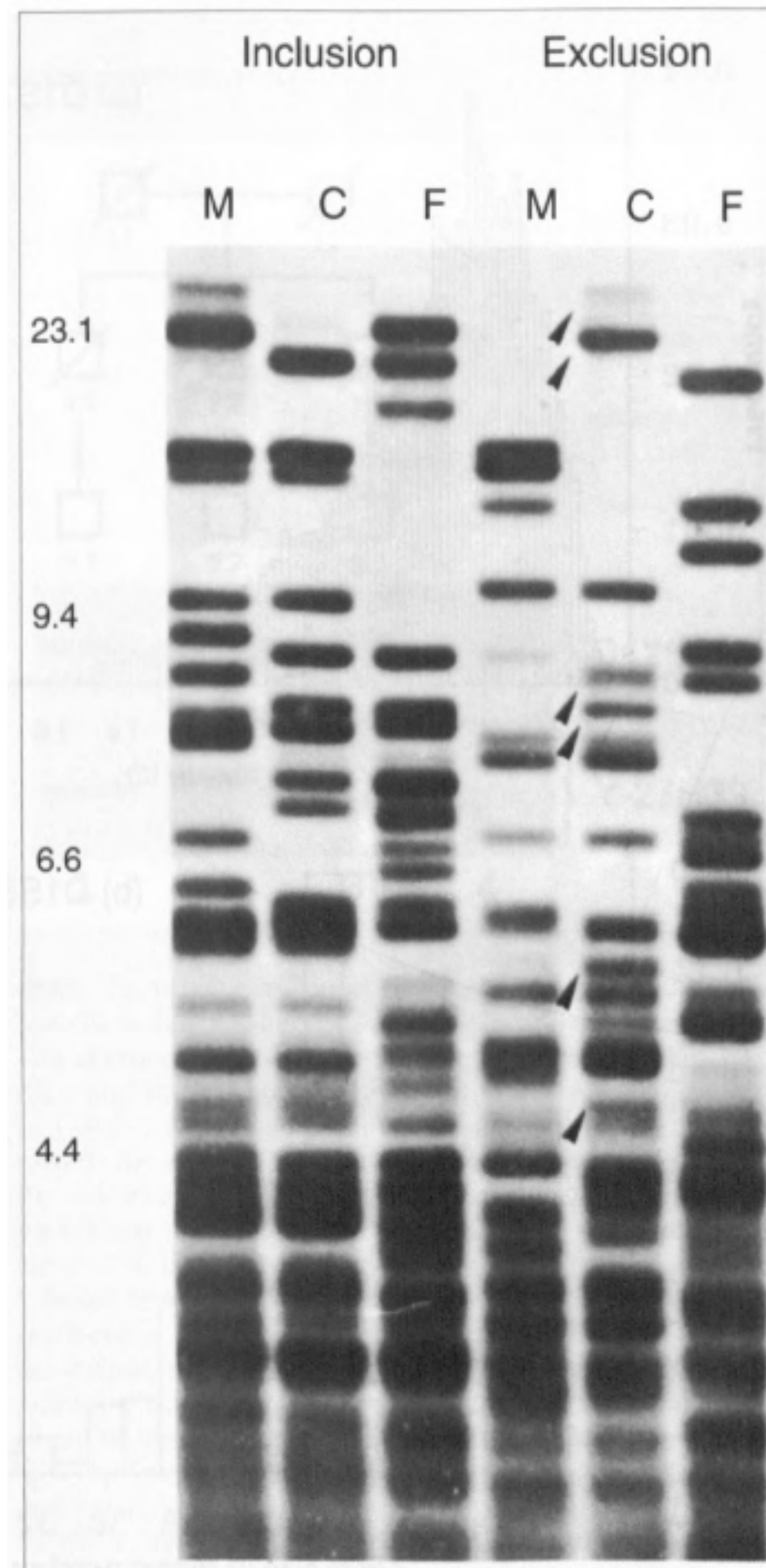
Copy number: number of times the unit is repeated

**STRs have mutation rates up to 10'000
times higher than point mutations**

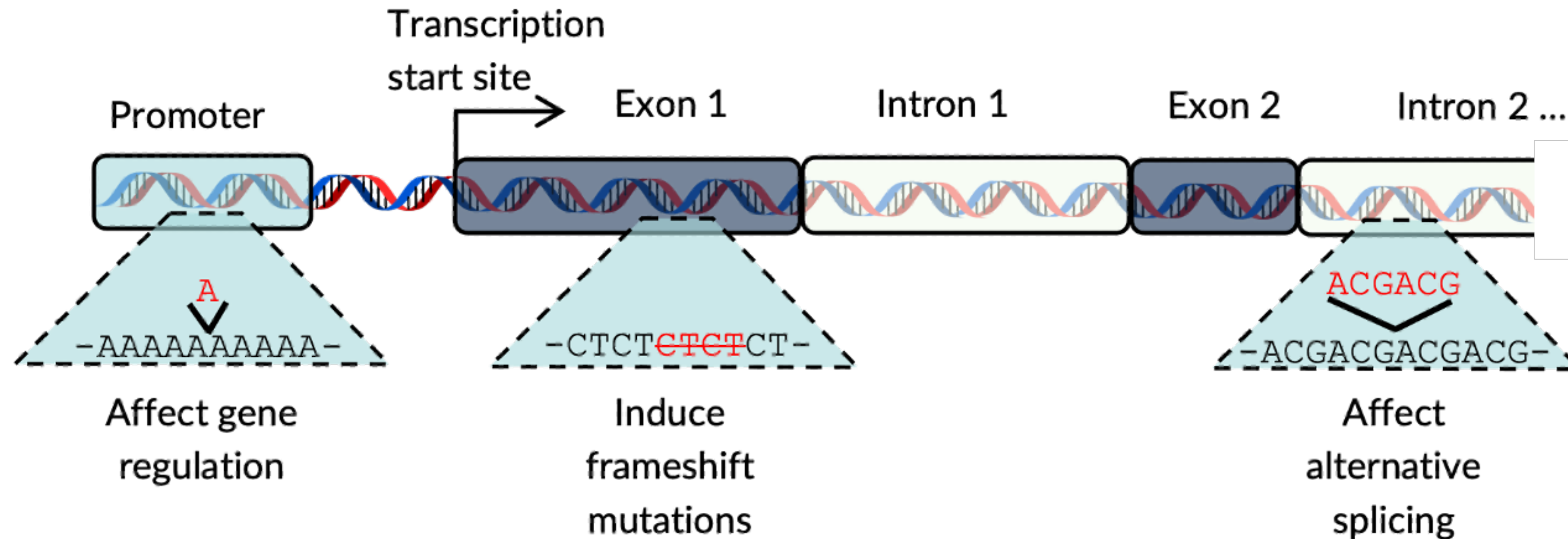
Slipped-strand mispairing



Classic use case: paternity testing



Functional consequences of STR mutations



We need special tools to genotype STRs!

e.g. Table 1 in <https://doi.org/10.1016/j.gde.2017.01.012>

- We will use GangSTR today: <https://github.com/gymreklab/GangSTR>

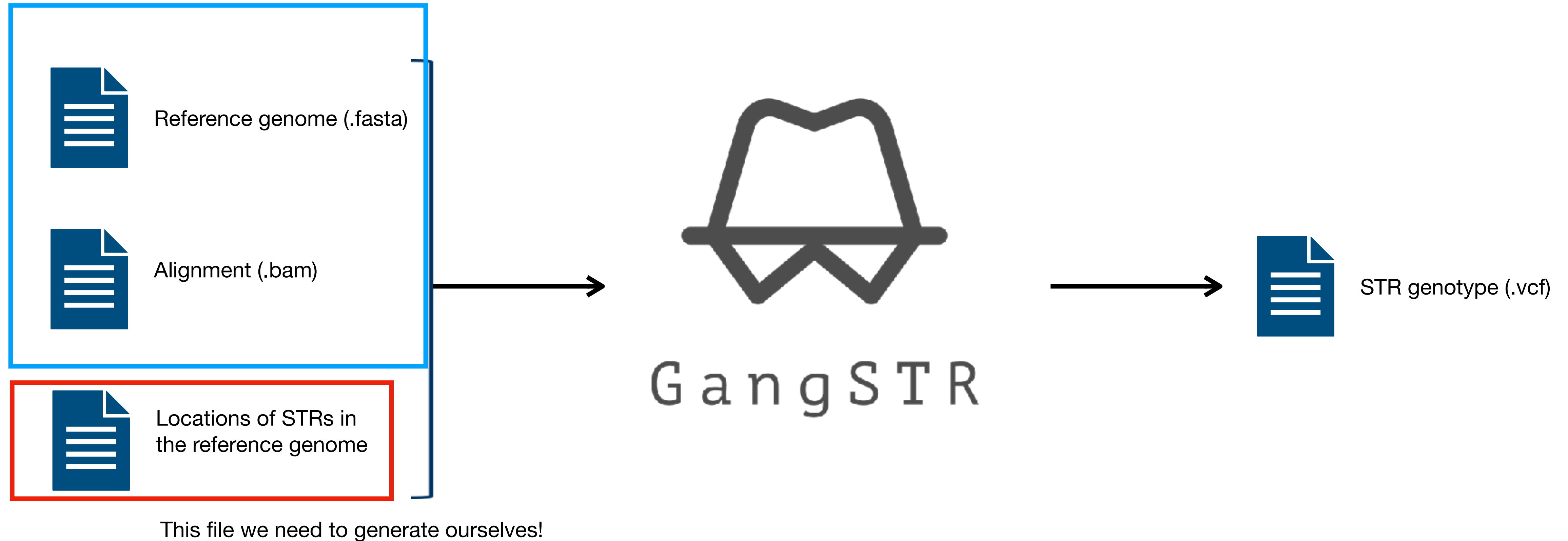


G a n g S T R

Mousavi, Nima, Sharona Shleizer-Burko, Richard Yanicky, and Melissa Gymrek. 'Profiling the Genome-Wide Landscape of Tandem Repeat Expansions'. *Nucleic Acids Research* 47, no. 15 (5 September 2019): e90–e90. <https://doi.org/10.1093/NAR/GKZ501>.

GangSTR intput/output

These files we have!
Reference: available online
Alignment: output of sequencing experiment



Detecting (short) tandem repeats in genomic sequence

- Computationally expensive for large sequences (especially imperfect repeats)
- Many different algorithms exist, each with different strengths/weaknesses

...GCTACGTACCTACTAACCTACCTACCTAA...

TR unit alignment

ACGTACCT
AC-TACCT
ACCTACCT

Tandem repeat annotation library (TRAL)

- A python library aimed at addressing challenges in repeat detection:
 - Allows for the running of several repeat detection algorithms on the input sequence
 - Calculates scores for each detected repeat for filtering
 - Can make the set of detected repeats non-redundant

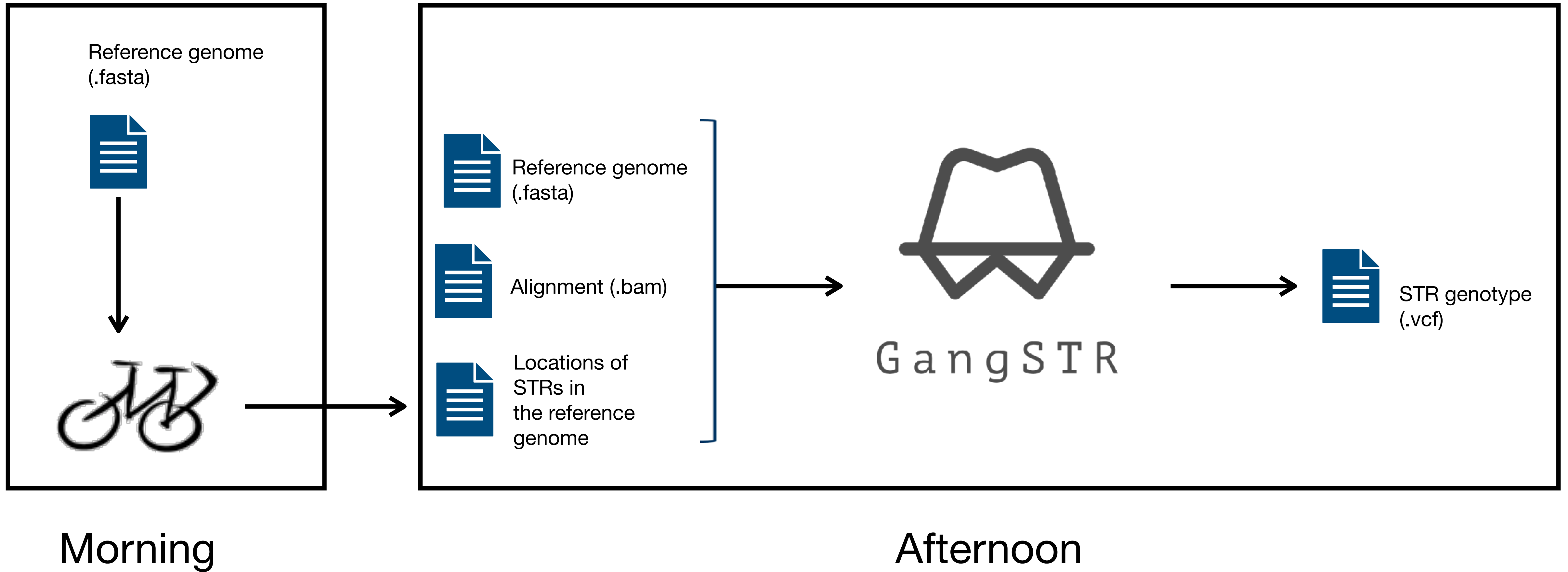


Schaper, Elke, Alexander Korsunsky, Jūlija Pečerska, Antonio Messina, Riccardo Murri, Heinz Stockinger, Stefan Zoller, Ioannis Xenarios, and Maria Anisimova. 'TRAL: Tandem Repeat Annotation Library'. *Bioinformatics* 31, no. 18 (15 September 2015): 3051–53. <https://doi.org/10.1093/bioinformatics/btv306>.

Summary

- Short tandem repeats can mutate rapidly, which can have functional consequences
- Repeat *genotyping*: determining how long repeats are in a given sequencing sample (e.g. using GangSTR)
- Repeat *detection*: determining where repeats are located and how long they are in the reference sequence (e.g. using TRAL)

Workflow for today



Task 1

STR background reading

Gymrek, Melissa. 'A Genomic View of Short Tandem Repeats'.
Current Opinion in Genetics and Development 44 (1 June 2017):
9–16. <https://doi.org/10.1016/j.gde.2017.01.012>.

- Read the following sections of 'A genomic view of short tandem repeats':
 - Abstract + Introduction
 - Genotyping STRs from high-throughput sequencing data
- Afterwards, you should be able to answer the following questions:
 - Why is STR variation relevant to health and disease?
 - What are some of the challenges in analysing STRs from NGS data?

Task 2

Introduction to TRAL

Schaper, Elke, Alexander Korsunsky, Jūlija Pečerska, Antonio Messina, Riccardo Murri, Heinz Stockinger, Stefan Zoller, Ioannis Xenarios, and Maria Anisimova. 'TRAL: Tandem Repeat Annotation Library'. *Bioinformatics* 31, no. 18 (15 September 2015): 3051–53. <https://doi.org/10.1093/bioinformatics/btv306>.

- Read 'TRAL: tandem repeat annotation library'
- Afterwards, you should be able to answer the following questions:
 - Why should you use multiple tandem repeat detection algorithms to look for repeats in biological sequence?
 - What different functionalities does TRAL provide?

Task 3

STR detection in the *APC* gene

- Set up the 'BIO392_STRs' conda environment and activate it (see appendix 1)
- Start a jupyter server and open 'scripts/BIO392_TRAL.ipynb' (see appendix 2)
- Follow along with the steps in BIO392_TRAL.ipynb

Appendix 1

Setting up the BIO392_STRs conda env

1: Make sure you are in the right directory

2: type 'ls + ENTER' you should see this file

```
maxverbiest@clt-mob-n-2720 2022-09-30 % ls
data          environment.yaml  results        scripts
maxverbiest@clt-mob-n-2720 2022-09-30 % conda env create -f environment.yaml
Collecting package metadata (repodata.json): done
Solving environment: done

==> WARNING: A newer version of conda exists. <==
  current version: 4.10.3
  latest version: 22.9.0

Please update conda by running

    $ conda update -n base -c defaults conda

Downloading and Extracting Packages
openjpeg-2.5.0           528 KB | ##### | 100%
ca-certificates-2022     150 KB | ##### | 100%
libcxx-1.17              1.3 MB | ##### | 100%
```

...

```
Installing collected packages: contigobj, biopython, tral
Successfully installed biopython-1.79 contigobj-5.0.6 tral-2.0
done
#
# To activate this environment, use
#
#     $ conda activate BIO392_STRs
#
# To deactivate an active environment, use
#
#     $ conda deactivate
#
maxverbiest@clt-mob-n-2720 2022-09-30 % conda activate BIO392_STRs
(BIO392_STRs) maxverbiest@clt-mob-n-2720 2022-09-30 %
```

4: The environment is created and ready to use!
type 'conda activate BIO392_STRs + ENTER' to activate it

Appendix 2

Starting a jupyter server

1: Make sure you are in the right directory

2: type 'jupyter notebook + ENTER'

```
(BI0392_STRs) maxverbiest@clt-mob-n-2720 2022-09-30 % jupyter notebook
[I 13:55:35.924 NotebookApp] Serving notebooks from local directory: /Users/maxverbiest/PhD/projects/BI0392/2022-09-30
[I 13:55:35.924 NotebookApp] Jupyter Notebook 6.3.0 is running at:
[I 13:55:35.924 NotebookApp] http://localhost:8888/
[I 13:55:35.924 NotebookApp] Use Control-C to stop this server and shut down all kernels (twice to skip confirmation).
/Users/maxverbiest/miniconda3/envs/BI0392_STRs/lib/python3.6/json/encoder.py:257: UserWarning: date_default is deprecated since jupyter_client 7.0.0. Use j
upyter_client.jsonutil.json_default.
  return _iterencode(o, 0)
```

3: The jupyter interface should open up in your browser
If not, you can copy this URL and go there yourself